

PROMPT SPACE OPTIMIZING FEW-SHOT REASONING SUCCESS WITH LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Prompt engineering is an essential technique for enhancing the abilities of large language models (LLMs) by providing explicit and specific instructions. It enables LLMs to excel in various tasks, such as arithmetic reasoning, question answering, summarization, relation extraction, machine translation, and sentiment analysis. Researchers have been actively exploring different prompt engineering strategies, such as Chain of Thought (CoT), Zero-CoT, and In-context learning. However, an unresolved problem arises from the fact that current approaches lack a solid mathematical solution for determining optimal prompts. To address this issue in prompt engineering, we propose a new and effective approach called Prompt Space. Our methodology utilizes text embeddings to obtain basis vectors by matrix decomposition, and then constructs a space for representing all prompts. Prompt Space significantly outperforms state-of-the-art prompt paradigms on ten public reasoning benchmarks. Notably, without the help of the CoT method and the prompt "Let's think step by step", Prompt Space shows superior performance over the few-shot method. Overall, our approach provides a robust and effective mathematical framework for selecting simple and effective prompts. This advancement marks a significant step towards improving prompt engineering for a wide variety of applications in LLMs.

1 INTRODUCTION

Prompt engineering becomes a relatively new and hot discipline for designing and optimizing prompts to effectively use large language models (LLMs) for a wide variety of applications and research domains (Brown et al., 2020; Thoppilan et al., 2022; Zhou et al., 2022; Sun et al., 2022; Dong et al., 2022). Researchers explore the use of simple and specific instructions to enhance the performance of LLMs on complex tasks, including arithmetic and commonsense reasoning, as well as question answering (Chowdhery et al., 2022; Scao et al., 2022; Ouyang et al., 2022; Bai et al., 2022). Developers strive to design robust and effective prompts either manually (Schick & Schütze, 2020; Reynolds & McDonell, 2021) or automatically (Gao et al., 2020) that interface with LLMs and other tools (Wu et al., 2023; Xie et al., 2023). The goal is to uncover the full potential of LLMs across various domains, enabling them to tackle complex tasks with improved performance and accuracy.

To elicit the reasoning ability of LLMs, Wei et al. (2022) has proposed the concept of the chain-of-thought (CoT) prompting. Unlike traditional input-output exemplars, the CoT prompting creates a series of intermediate reasoning steps that guide LLMs through a complex problem. This approach enables LLMs to develop a reasoning path that decomposes the complex problem into multiple reasoning steps. Notably, the CoT prompting demonstrates that the reasoning ability of LLMs perfectly matches the scaling laws, with the reasoning ability of LLMs significantly increasing with the size of the PaLM 540B model. Inspired by the CoT prompting, several works explore methods to enhance LLMs' reasoning abilities with simple techniques. Kojima et al. (2022) introduces the "Let's think step by step" prompt, which helps LLMs adopt a step-by-step thinking approach, leading to the final answer. Their approach, known as Zero-shot-CoT, successfully generates a reasoning path in zero-shot reasoning scenarios. In practice, the CoT prompting has showed better performance than Zero-shot-CoT (Wei et al., 2022; Kojima et al., 2022). However, the CoT prompting involves huge efforts in the manual design of both questions and related reasoning chains. To avoid the manual method, Zhang et al. (2022) proposes an automatic CoT prompting, called Auto-CoT. It applies a clustering algorithm to identify representative questions for each cluster and generates reasoning

chains using the Zero-shot-CoT method for each question. By combining the prompt, "let's think not just step by step, but also one by one," Auto-CoT explores optimal prompts across various reasoning tasks.

The previous works on CoT have greatly contributed to our understanding of effective prompts for improving the reasoning ability of LLMs. However, these works have certain limitations, such as the lack of guidance on finding optimal prompts for reasoning tasks. In this paper, we propose a novel approach called Prompt Space that overcomes these limitations and simultaneously leverages the strengths of previous works. Our approach starts by embedding questions and then utilizes matrix decomposition to yield basis vectors, or basis questions. These basis questions are used to construct a space that can represent all questions. With Zero-shot-CoT, we combine these basis questions with every question to automatically generate reasoning demonstrations for LLMs. Our approach offers a promising solution to find optimal prompts on reasoning tasks and significantly improves the few-shot reasoning of LLMs.

Prompt Space surpasses the performance of current prompt paradigms on ten public reasoning benchmarks. Our work uncovers critical insights into the impact of the number of basis questions on reasoning tasks. Additionally, we identify the relationship between the selected questions and the reasoning ability of LLMs, and investigate how to determine the optimal number of exemplars for each reasoning task. Extensive experiments demonstrate that our approach establishes a reliable and mathematical methodology for selecting simple and effective prompts. Our goal is to not only design the robust and effective prompts for challenging reasoning tasks, but also highlight the significance of carefully exploring and analyzing the optimal prompts for unlocking the potential of LLMs in a wide variety of applications.

2 RELATED WORK

2.1 CHAIN-OF-THOUGHT PROMPTING

Chain-of-thought (CoT) prompting is an effective method to elicit the reasoning ability of LLMs through a chain of thought, where a series of intermediate reasoning steps are used to generate the answer (Wei et al., 2022). This approach has been shown to significantly improve the performance of LLMs on complex reasoning tasks. To further enhance their performance, self-consistency (SC) has been introduced, which replaces the standard greedy decoding of the LLM output with a stochastic output space ensemble Wang et al. (2022b). Existing studies on the CoT prompting can be broadly divided into two categories: manually constructed the CoT prompting and automatically generated the CoT prompting. Our work aims at providing a robust and mathematical framework for selecting simple and effective prompts.

2.2 AUTOMATICALLY GENERATED CoT PROMPTS

For enhancing CoT reasoning in LLMs, several previous works have explored the idea of self-generating a chain of thought Kojima et al. (2022); Zhang et al. (2022); Zhou et al. (2022); Hebenstreit et al. (2023). Kojima et al. (2022) finds that using specific phrases, like "Let's think step by step", as a prompt can guide LLMs to generate reasoning steps without any few-shot hand-crafting exemplars. Following this work, Zhou et al. (2022) proposes a framework called Automatic Prompt Engineer (APE) for generating and selecting instructions automatically. APE addresses the instruction generation problem by using LLMs to generate and search for candidate solutions.

Additionally, some studies implement Zero-shot-CoT to generate the reasoning process in their demonstration (Kojima et al., 2022). Zhang et al. (2022) proposes a novel method called Auto-CoT, for automatically creating the Chain of Thought (CoT) prompting in LLMs. This method samples diverse questions and reasoning chains to construct effective demonstrations for LLMs. It can elicit chain-of-thought reasoning without impairing performance and eliminating the need for hand-crafting prompts. In contrast, Shao et al. (2023) utilizes seed demonstrations to automatically synthesize more examples through forward and backward processes. Inspired by these works, we construct a space with text embeddings and the matrix decomposition to represent all questions. We also utilize Zero-shot-CoT to generate chains of thought for prompt examples (Kojima et al., 2022).

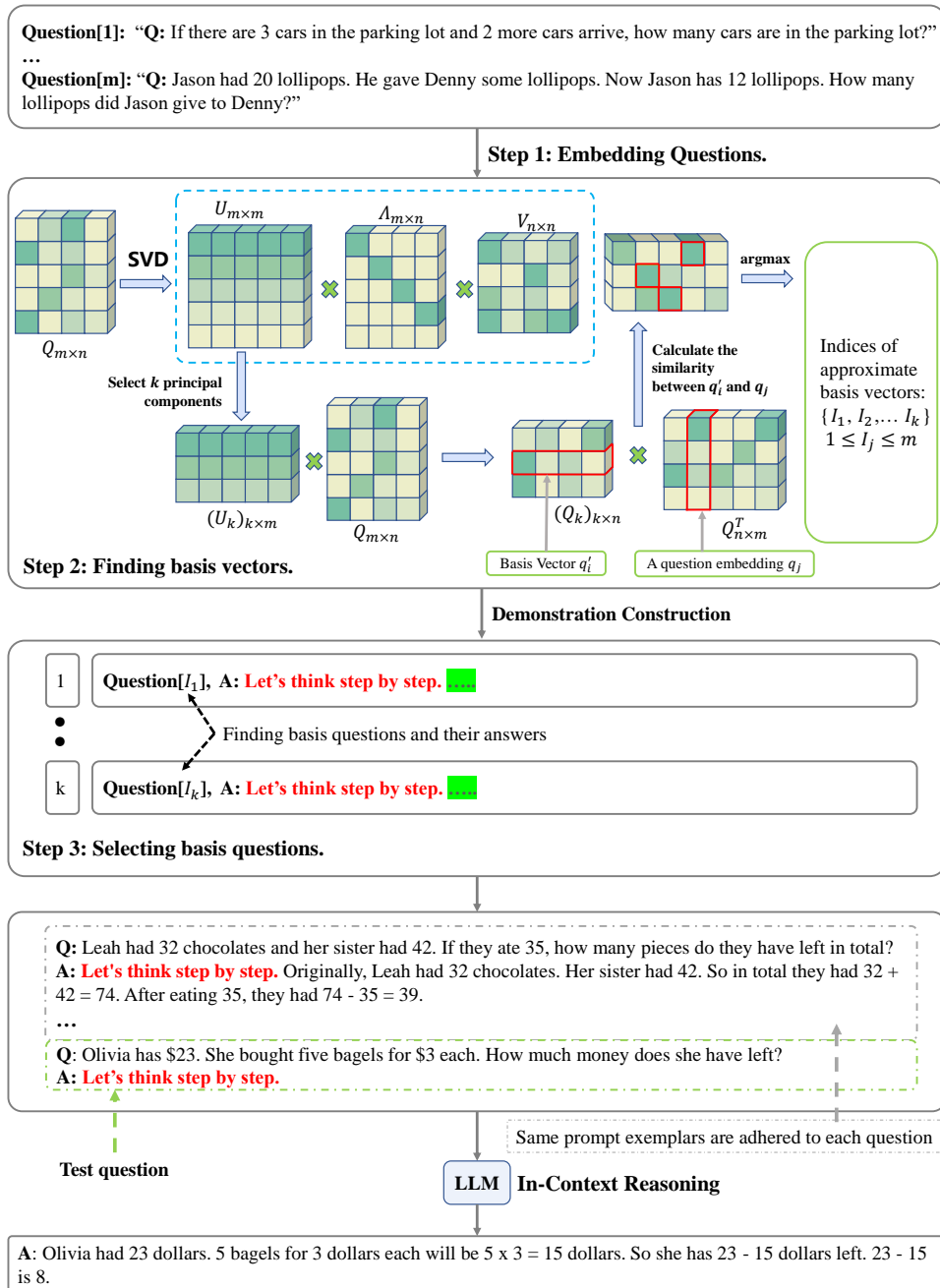


Figure 1: The schematic of Prompt Space. Prompt Space consists of three steps, including embedding questions, finding basis vectors, and selecting basis questions. First, all questions in a reasoning dataset are encoded into the question matrix $Q_{m \times n}$. Second, the k basis vectors in the question matrix $Q_{m \times n}$ are calculated by SVD and PCA. Finally, the top k questions are selected to form the question space. The prompt exemplar is made up of the selected k questions and the test question. LLMs could follow the prompt exemplar and then develop a chain of reasoning steps to get the final answer.

2.3 EXAMPLE SELECTION

For designing prompts, certain studies demonstrate that the performance of LLMs is influenced by various factors, such as tasks, prompts, and model structures (Zhao et al., 2021; Lu et al., 2021; Su et al., 2022; Griffin et al., 2023; Jiang et al., 2022). The main challenge is to develop selection criteria that are both effective and generalizable based on empirical experiments. Wang et al. (2022b) shows that the sequence of reasoning steps is critical for achieving optimal performance. Additionally, Rubin et al. (2021) proposes a similarity-based selection method, which retrieves the most similar training instances as a prompt for a given test case. Furthermore, another approach proposed by Fu et al. (2022) suggests the selected prompts with more steps can significantly improve performance in the reasoning process. However, our method, Prompt Space, explores question embeddings to obtain basis questions in a reasoning task, which dramatically avoids using ineffective questions as a demonstration. This approach provides an innovative mathematical solution for selecting effective prompts, yielding more generalizable and comprehensive reasoning chains. Our Prompt Space aims to develop a deep understanding of how to design the CoT prompting.

3 PROMPT SPACE

In this work, we propose a novel method called Prompt Space, which automatically creates demonstrations with questions and reasoning chains. Prompt Space seeks to design an appropriate space for identifying the basis questions for building prompt exemplars. For a vector space V , its vector basis is defined as a subset $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ in V . These basis vectors are linearly independent in span V . Consequently, if $(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n)$ is a list of vectors in V , then these vectors form a vector basis if and only if every $x \in V$ can be uniquely written as

$$\mathbf{x} = c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_n\mathbf{v}_n, \quad (1)$$

where c_1, c_2, \dots, c_n are elements of the base field.

In Prompt Space, one vector represents a question embedding. By combining the basis questions with the test question, we create a demonstration that enables LLMs to effectively generate a chain of thought. Next, we will show how to select such basis questions for constructing Prompt Space.

Selecting effective prompts as an exemplar can significantly enhance the reasoning abilities of LLMs. To solve arithmetic problems, humans tend to learn from previous question-answer pairs and generalize them to solve similar problems. Inspired by this thought mechanism, our work aims to select more representative questions as an exemplar to facilitate LLMs in developing a chain of reasoning steps. We assume that there exists a real prompt space P with k dimensional vectors, where the selected representative questions can serve as basis vectors of this space. These basis vectors provide an effective solution for LLMs to reason through the problem space.

Principal component analysis (PCA) is a widely used algorithm for identifying the key components of extensive data features by geometric steps (Abdi & Williams, 2010). The implementation of PCA can efficiently compress a n -dimension matrix into a k -dimension matrix and obtain k principal vectors from the original space. Inspired by PCA, Prompt Space follows the below steps:

1. **Embedding Questions.** The question set of a task is $\mathcal{Q} = \{q_1, q_2, \dots, q_m\}$, where m is the number of questions in a task. The MiniLM-L6-v2 model f_{MiniLM} (Wang et al., 2020) encodes these questions as follows: $\mathbf{q}_i = f_{MiniLM}(q_i) \in \mathbb{R}^n$, for $i = 1, 2, \dots, m$. After the encoding process, the question matrix Q is created by putting together all question embeddings, i.e., $Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_m]^T \in \mathbb{R}^{m \times n}$.

Assuming the dimension (rank) of the prompt space P is k , the process of finding k basis vectors is identical to searching for the k principal components (questions) of the question matrix Q .

2. **Finding basis vectors.** We use Singular Value Decomposition (SVD) to calculate k basis vectors in the prompt space P (Wall et al., 2003). Using SVD, Q can be calculated as:

$$Q = U\Lambda V^T, \quad (2)$$

where U is denoted as a left singular matrix, $U = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m]^T \in \mathbb{R}^{m \times m}$, and $\mathbf{u}_i \in \mathbb{R}^{1 \times m}$ is the eigenvector of $QQ^T \in \mathbb{R}^{m \times m}$ (for $i = 1, \dots, m$). A complete proof is

Prompt Demonstrations on Last Letter	
Random Selection	Manual-CoT
<p>Q: Take the last letters of each words in "Amy Ella Amilcar Roman" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Amy" is "y". The last letter of "Ella" is "a". The last letter of "Amilcar" is "r". The last letter of "Roman" is "n". So, we have "y" + "a" + "r" + "n" = "yarn". Therefore, the answer is "yarn".</p> <p>Q: Take the last letters of each words in "Lara Greg Ethan Terrence" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Lara" is "a". The last letter of "Greg" is "g". The last letter of "Ethan" is "n". The last letter of "Terrence" is "e". Therefore, the answer is "agne".</p> <p>Q: Take the last letters of each words in "Luisa Tatiana Drake Hillary" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Luisa" is "a". The last letter of "Tatiana" is "a". The last letter of "Drake" is "e". The last letter of "Hillary" is "y". Therefore, the answer is "aaey".</p> <p>Q: Take the last letters of each words in "Gus Brock Ava Jenna" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Gus" is "s". The last letter of "Brock" is "k". The last letter of "Ava" is "a". The last letter of "Jenna" is "a". Therefore, the answer is "skaa".</p>	<p>Q: Take the last letters of the words in "Elon Musk" and concatenate them.</p> <p>A: The last letter of "Elon" is "n". The last letter of "Musk" is "k". Concatenating them is "nk". The answer is nk.</p> <p>Q: Take the last letters of the words in "Larry Page" and concatenate them.</p> <p>A: The last letter of "Larry" is "y". The last letter of "Page" is "e". Concatenating them is "ye". The answer is ye.</p> <p>Q: Take the last letters of the words in "Sergey Brin" and concatenate them.</p> <p>A: The last letter of "Sergey" is "y". The last letter of "Brin" is "n". Concatenating them is "yn". The answer is yn.</p> <p>Q: Take the last letters of the words in "Bill Gates" and concatenate them.</p> <p>A: The last letter of "Bill" is "l". The last letter of "Gates" is "s". Concatenating them is "ls". The answer is ls.</p>
Auto-CoT	Prompt Space
<p>Q: Take the last letters of each words in "Tim Candace Cecil Misael" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Tim" is "m". The last letter of "Candace" is "e". The last letter of "Cecil" is "l". The last letter of "Misael" is "l". So the final answer would be "mel".</p> <p>Q: Take the last letters of each words in "Alina Alessandra Amina Bianca" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Alina" is "a". The last letter of "Alessandra" is "a". The last letter of "Amina" is "a". The last letter of "Bianca" is "a". So the final answer is "aaaa".</p> <p>Q: Take the last letters of each words in "Felipe Heidi Nino Bradley" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Felipe" is "e". The last letter of "Heidi" is "i". The last letter of "Nino" is "o". The last letter of "Bradley" is "y". So, the final answer would be "eioly".</p> <p>Q: Take the last letters of each words in "Lacey Nora Debra Ashleigh" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Lacey" is "y". The last letter of "Nora" is "a". The last letter of "Debra" is "a". The last letter of "Ashleigh" is "h". So the final answer is "yahah".</p>	<p>Q: Take the last letters of each words in "Ever Gio Ella Ramesh" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Ever" is "r". The last letter of "Gio" is "o". The last letter of "Ella" is "a". The last letter of "Ramesh" is "h". So, the final answer is "roah".</p> <p>Q: Take the last letters of each words in "Marian Joanne Darrin Rohit" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Marian" is "n". The last letter of "Joanne" is "e". The last letter of "Darrin" is "n". The last letter of "Rohit" is "t". So the final answer would be "nent".</p> <p>Q: Take the last letters of each words in "Craig Dillon Troy Griselda" and concatenate them.</p> <p>A: Let's think step by step. The last letter of the first word is "g". The last letter of the second word is "y". The last letter of the third word is "a". The last letter of the fourth word is "a". Therefore, the final answer is "gyaa".</p> <p>Q: Take the last letters of each words in "Liliana Quincy Bart Makayla" and concatenate them.</p> <p>A: Let's think step by step. The last letter of "Liliana" is "a". The last letter of "Quincy" is "y". The last letter of "Bart" is "t". The last letter of "Makayla" is "a". So the final answer is "ayta".</p>

Figure 2: The demonstrations of Last Letter on difference methods, including Random selection, Manual-CoT, Auto-CoT, and Our Prompt Space.

shown in Appendix A. Similarly, V is the right singular matrix, which can be written as $V = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m]^T \in \mathbb{R}^{n \times n}$, $\mathbf{v}_i \in \mathbb{R}^{n \times 1}$ (for $i = 1, \dots, n$) is the eigenvector of QQ^T . Next, the k principal components of Q can be obtained:

$$Q_k = U_k Q, \quad (3)$$

where $U_k = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k]^T \in \mathbb{R}^{k \times m}$, and $Q_k \in \mathbb{R}^{k \times n}$ contains the top k principal components ranked by related eigenvalues. The row vectors in $Q_k = [\mathbf{q}'_1, \mathbf{q}'_2, \dots, \mathbf{q}'_k]^T \in \mathbb{R}^{k \times n}$ are k basis vectors of the prompt space P .

3. Selecting basis questions. In this step, we select the top k questions from the question matrix Q as basis questions, whose embeddings are closest to these basis vectors. It can result in:

$$f(x) = \operatorname{argmax}(x \cdot Q^T), \text{ for } x \in \{\mathbf{q}'_i \in \mathbb{R}^{1 \times n}, i = 1, 2, \dots, k\} \quad (4)$$

where $\operatorname{argmax}(\bullet)$ is to calculate the maximum similarities between question embeddings and basis vectors (i.e. cosine similarity) (Sidorov et al., 2014). Finally, we can generate the prompt exemplar, including k basis questions and the original question in Q , to get the final output (answer).

Figure 1 shows that an example of Prompt Space produces basis questions to solve an arithmetic problem. By following the three steps, we could select k basis questions and then combine them with the text question. To assist LLMs in generating the final output, we also use the prompt, "Let's think step by step". Throughout the process, we still opt to automatically create the prompt rather than manually design it. As a result, LLMs can generate a step-by-step thought process for arriving at the answer. The algorithm of Prompt Space is shown in Appendix B.

Prompt Space has several attractive properties as an approach for enhancing reasoning in LLMs.

1. Prompt Space enables LLMs to identify optimal prompts for a range of reasoning tasks and efficiently generate final outputs.
2. Prompt Space provides a robust mathematical framework for designing the prompt. It can suggest the optimal number of exemplars to improve the reasoning abilities of LLMs. Our method provides valuable insights into effective prompting strategies for achieving successful outcomes.
3. Prompt Space has the potential to be utilized in a variety of few-shot tasks through prompt engineering, including but not limited to translation, summarization, and expansion.

4 EXPERIMENTS

We briefly describe the experimental setup and highlight the main results. Additional experimental details and results can be found in Appendices C and D.

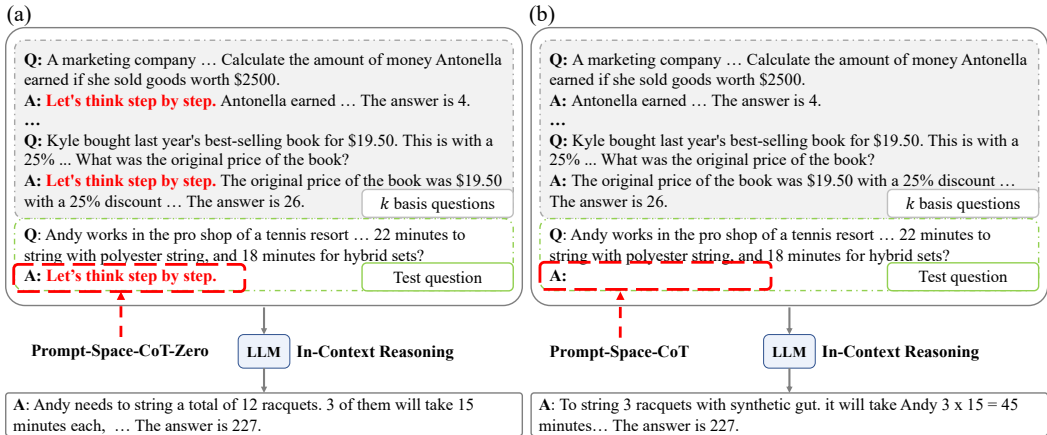


Figure 3: Prompt-Space-CoT-Zero (including the prompt, “Let’s think step by step” prompt) and Prompt-Space-CoT (not including the prompt, “Let’s think step by step”) with an input-output exemplar of an LLM.

Prompt Space is evaluated on three categories of reasoning tasks, namely arithmetic reasoning, commonsense reasoning, and symbolic reasoning. The experiment demonstrates Prompt Space on various tasks: 1. Prompt Space outperforms the state-of-the-art baselines on these tasks. 2. Prompt Space can efficiently construct a space and find its basis questions for each task. 3. Prompt Space can determine the optimal number of basis questions for significantly improving the performance of LLMs on each dataset. 4. Prompt Space relies on the selection of embedding models.

4.1 EXPERIMENTAL SETUP

Tasks and Datasets. Prompt Space is studied on ten standard datasets from three categories of reasoning tasks:

1. Arithmetic reasoning contains six datasets: (1) AddSub (Hosseini et al., 2014), (2) MultiArith (Roy et al., 2015), (3) SingleEq (Koncel-Kedziorski et al., 2015), (4) AQUA-RAT (Ling et al., 2017), (5) SVAMP (Patel et al., 2021), (6) GSM8K (Cobbe et al., 2021). These datasets are sorted by release time. SingleEq and AddSub have plenty of easier problems, while MultiArith, AQUA-RAT, SVAMP, and GSM8K are more difficult and require multi-step reasoning steps.

2. Commonsense reasoning: (1) CommonsenseQA (CSQA) (Talmor et al., 2019), (2) StrategyQA (STQA) (Geva et al., 2021). CSQA is a challenging dataset for commonsense question answering. Its questions contain complex semantics that often requires prior knowledge. STQA requires multi-step reasoning with an inferred strategy in the question.

3. Symbolic reasoning: (1) Last Letter Concatenation (Letter) and (2) Coin Flip (Coin) (Wei et al., 2022). Last Letter Concatenation asks the model to concatenate the last letters of words in a name. We generate full names by randomly concatenating names from samples. Coin Flip requires the model to answer whether a coin is still heads up after people either flip it or do not flip it. In this work, we consider an *out-of-domain* test set, where examples had more steps than those in exemplars.

The detailed description of each dataset is shown in Appendix C.1.

Baselines. We compare our Prompt Space with five baseline methods: Few-shot (Wei et al., 2022), Manual-CoT (Wei et al., 2022), Zero-shot (Kojima et al., 2022), Zero-shot-CoT (Kojima et al., 2022), and Auto-CoT (Zhang et al., 2022). Few-shot easily selects question-answer pairs as a demonstration for feeding to LLMs. Manual-CoT involves manually creating a series of reasoning chains as a demonstration to elicit the reasoning ability of LLMs. Zero-shot is a standard prompting technique for evaluating the abilities of LLMs. Zero-shot-CoT randomly selects questions as demonstrations, and then uses the prompt "Let’s think step by step". Additionally, Auto-CoT utilizes clustering techniques to sample questions and generate demonstrations with the Zero-shot-CoT method. To ensure fair comparisons with the baselines, we run experiments with consistent in-context exemplars and a constant seed across all methods and datasets. Few-shot and Manual-CoT select the examples

by human, while Auto-CoT select the examples by the K-means clustering algorithm. Our Prompt Space uses the same rationales with Zero-CoT and Auto-CoT not Manual-CoT. Our CoT is generated by LLMs not humans. Figure 2 shows demonstrations of CSQA on difference methods including Random selection, Manual-CoT, Auto-CoT and Our Prompt Space. Please refer to Appendix C.2 for detailed baselines.

Implementation. We use the gpt-35-turbo version of the public ChatGPT model from the OpenAI API with 175 billion parameters (Brown et al., 2020; Ouyang et al., 2022). We select this LLM because it has better performance than the text-davinci-002 version of GPT-3, as reported in (OpenAI, 2023; Bai et al., 2022). In the decoding process, we set the temperature to 0 and use a greedy searching algorithm to obtain results. For zero-shot approaches, our results are deterministic. Following Wei et al. (2022), we set the number of demonstrations k to 8, except for AQUA-RAT (4) and Letter (4), StrategyQA (6), and CSQA (7). However, Our Prompt Space can determine the optimal number of basis questions for each task. In the following sections, we will present a detailed analysis and provide further insights into the selection of basis questions. The selected embedding models are T5 models (base/large/XL/XXL) (Raffel et al., 2020), E5 models (small/base/large) (Wang et al., 2022a) and MiniLM-L6-v2 model (Wang et al., 2020). The embedding size of each question in all T5 models is 768, while for E5 models (small, base, large), their embedding sizes are 384, 768, 1024 respectively. Our MiniLM-L6-v2 model encodes questions with an embedding size of 384. Please refer to Appendix C.3 for detailed model descriptions. In our approaches, we investigate two types of Prompt Space shown in Fig. 3. The first type combines CoT with the ‘‘Let’s think step by step’’ prompt, denoted as Prompt-Space-CoT-Zero. In contrast, the second type only uses CoT, namely Prompt-Space-CoT.

4.2 MAIN RESULTS

In the experiments, we evaluate Prompt Space on ten datasets from three categories of reasoning tasks. Due to the greedy decoding, the main results show deterministic results without error bars. Notably, Table 1 and 2 show that Prompt Space achieves superior performance over the state-of-the-art methods on ten reasoning tasks, respectively. Compared to Auto-CoT, Prompt space with the optimal number of exemplars achieves up to average 3.2% in Table 2.

Prompt Space vs Few-shot. Table 1 summarizes comparisons between our approach (Prompt Space) and two baselines (Zero-shot and Few-shot) for each dataset. In Table 1, Prompt Space doesn’t include CoT and the ‘‘Let’s think step by step’’ prompt, and just selects basis questions as a demonstration. Our results show that Prompt Space with the same settings achieves up to average 2.3%, 2% over Zero-shot and Few-shot on ten reasoning datasets, respectively. Especially, Prompt Space, with the optimal number of exemplars, achieves up to average 3.3%, 3% over Zero-shot and Few-shot on ten reasoning datasets, respectively. The most significant improvement is observed in the STQA and Letter datasets, with a relative increase of 13.5%, 112.5% over Few-shot, respectively. Moreover, Prompt Space outperforms two baselines on eight out of ten reasoning datasets.

Table 1: Accuracy (%) comparison of Prompt Space with two baselines on ten reasoning datasets. Two baselines are Zero-shot and Few-shot, respectively. Ten benchmark datasets contain three categories, including arithmetic reasoning, commonsense reasoning, and symbolic reasoning. The last column shows average scores. See Appendix C for a detailed setup.

Model	Arithmetic						Commonsense		Symbolic		Avg
	AddSub	MultiArith	SingleEq	AQUA-RAT	SVAMP	GSM8K	CSQA	STQA	Letter	Coin	
Zero-shot	87.6	80.0	87.4	27.6	74.0	22.9	73.6	61.3	0.8	22.8	57.8
Few-shot	85.8	82.7	89.4	30.7	76.1	24.0	79.3	54.0	1.6	57.0	58.1
Prompt Space w/o CoT-Zero	89.4	83.7	88.8	32.7	75.2	25.2	79.1	61.3	3.4	62.0	60.1
Prompt Space w/o CoT-Zero (best)	89.9(10)	86.3 (9)	88.8 (8)	32.7 (3)	75.6 (6)	25.9 (6)	80.0 (8)	62.8 (10)	5.2 (9)	63.8 (6)	61.1

Arithmetic Reasoning. Our approach substantially outperforms the three baselines on five arithmetic reasoning tasks except for AddSub in Tab. 2. Importantly, our Prompt Space with the same settings achieves score gains of 1.8%, 1.2%, 2% and 2.1% over the previous state-of-the-art methods on MultiArith, SingleEq, SVAMP, and GSM8K, respectively. Although Prompt Space doesn’t show competitive performance on AddSub, it is close to Auto-CoT. Additionally, Prompt Space achieves the highest performance on AQUA-RAT, SVAMP, and GSM8K, indicating that it can solve more complex arithmetic reasoning. The difference between Prompt-Space-CoT-Zero and Prompt-Space-CoT is

Table 2: Accuracy (%) comparison of Prompt Space with four state-of-the-art methods on ten reasoning datasets. These datasets contain three categories, such as arithmetic reasoning, commonsense reasoning, and symbolic reasoning. The last column records average scores. See Appendix C for a detailed setup.

Model	Arithmetic					Commonsense		Symbolic		Avg	
	AddSub	MultiArith	SingleEq	AQUA-RAT	SVAMP	GSM8K	CSQA	STQA	Letter		Coin
Zero-shot-CoT	82.5	96.0	90.4	38.2	76.5	57.1	72.0	57.6	71.0	64.4	69.9
Manual-CoT	86.8	97.0	90.9	45.3	80.2	75.8	72.2	61.7	78.8	100	78.9
Auto-CoT	88.5	96.0	90.2	46.5	78.2	74.1	72.7	62.4	72.6	100	78.1
Prompt-Space-CoT-Zero	87.3	98.0	89.2	36.2	82.2	72.4	71.1	63.3	82.0	100	78.3
Prompt-Space-CoT	87.9	98.8	92.1	40.6	81.0	77.9	74.1	62.5	79.6	100	79.5
Prompt-Space-CoT (best)	87.9 (8)	98.8(8)	92.5(4)	48.8(7)	82.6(10)	77.9(8)	77.9(8)	64.4(9)	84.4(1)	100(3)	81.3

trivial, approximately 2%. Overall, the average score of Prompt-Space-CoT surpasses that of the three baselines on all arithmetic reasoning datasets, indicating its superior performance.

Commonsense Reasoning. Prompt Space significantly outperforms the prior state-of-the-art Auto-CoT over two commonsense reasoning datasets. Our approach with the same settings achieves respective improvements of 1.9%, 1.6% over Manual-CoT and 1.4%, 0.9% over Auto-CoT. Compared to Zero-shot, Zero-shot-CoT and Manual-CoT don’t elicit better commonsense reasoning, while Prompt Space leverages the CoT method to dramatically increase performance instead of decreasing it. These results demonstrate that Prompt Space can improve performance on commonsense reasoning tasks requiring prior knowledge.

Symbolic Reasoning. The performance of Prompt Space achieves a significant increase of 3.2% over Manual-CoT and 9.4%, over Auto-CoT on the Letter dataset, respectively. Interestingly, the accuracy of Manual-CoT, Auto-CoT, and our approach reaches to 100% on the Coin Flip dataset. The result indicates that Prompt Space dramatically enhances the reasoning abilities of LLMs on symbolic tasks.

4.3 EFFECT OF EMBEDDING MODELS

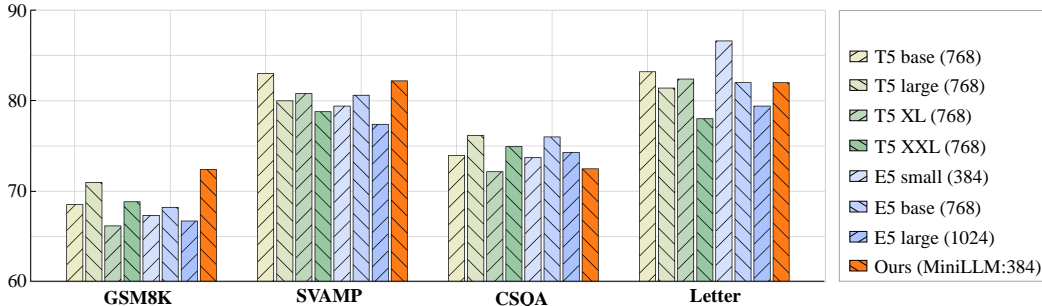


Figure 4: Accuracy (%) of Prompt-Space-CoT-Zero with different embedding models on three types of reasoning tasks: arithmetic reasoning (GSM8K, SVAMP), commonsense reasoning (CSQA), and symbolic reasoning (Letter). The embedding models are T5 (base/large/XL/XXL), E5 (small/base/large), and MiniLM-L6-v2 (ours), respectively.

Figure 4 shows that the increase of embedding size cannot improve the performance of Prompt Space on various reasoning tasks. Besides, the appropriate embedding size could be 768 in T5 and E5 models. As T5 models increase their model size, the performance of Prompt Space decreases significantly. Moreover, the solving rate of Prompt Space exhibits clear fluctuations on different embedding models.

4.4 FURTHER ANALYSIS OF BASIS QUESTIONS

Figure 5 illustrates the performance of Prompt Space with different basis questions on nine datasets. Our results reveal that the appropriate number of basis questions is 8 on arithmetic reasoning tasks

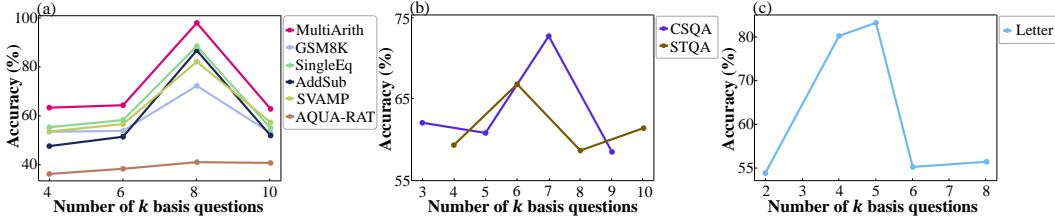


Figure 5: Accuracy (%) of Prompt-Space-CoT-Zero with various numbers of basis questions on nine reasoning datasets.

except for AQUA-RAT, while that of basis questions is approximately 6 or 7 on commonsense reasoning tasks. Interestingly, the AQUA-RAT and Letter datasets exhibit a preference of a smaller number of basis questions, which indicates their space could be spanned by just four or five basis vectors. Overall, our findings demonstrate that the appropriate number of basis questions can significantly improve performance, which indicates that there exist basis vectors (questions) in the prompt space. However, there remains a challenge that we cannot automatically determine the optimal number of basis questions for each dataset. More analysis about basis questions of Prompt-Space-CoT is shown in Appendix D.1. Besides, we provide more visualizations of Prompt Space and the constructed demonstrations in Appendix D.2 and E, respectively.

4.5 EFFECT OF QUESTION SEQUENCE

Table 3: Accuracy (%) of Prompt-Space-CoT-Zero with different question sequences on three types of reasoning tasks: arithmetic reasoning (GSM8K, SVAMP), commonsense reasoning (CSQA), and symbolic reasoning (Letter). Our model ranks the sequence of basis questions by their eigenvalue scores from largest to smallest. Three different sequences are evaluated: (a) original sequence, (b) reversed sequence, (c) random sequence.

Sequence	Arithmetic		Commonsense	Symbolic	Avg
	SVAMP	GSM8K	CSQA	Letter	
Random sequence	81.4	70.5	73.1	81.40	76.6
Reversed sequence	82.1	70.9	74.8	82.0	77.4
Original sequence (ours)	82.2	72.4	72.5	82.0	77.3

Table 3 shows that Prompt Space achieves better performance than other cases, when the basis questions are sorted in ascending order of their eigenvalues. However, the descending sort (original sequence) has superior performance over baselines on three out of four benchmarks. Furthermore, the difference between the original sequence and the reverse sequence is trivial ($\sim 0.1\%$). Thus, these findings suggest that the descending sort is an acceptable approach used in our experiments.

5 CONCLUSION

In this paper, we propose a novel prompting method, namely Prompt Space, to explore the selection of prompts for enhancing reasoning in LLMs. For any dataset, Prompt Space can map its questions onto a real space for determining basis questions as a demonstration. Through experiments on arithmetic, commonsense, and symbolic reasoning tasks, we find that the demonstrations constructed by Prompt Space can significantly improve the reasoning abilities of LLMs on ten public benchmarks. Furthermore, without the help of the CoT method and the "Let's think step by step" prompt, Prompt Space also exhibits superior performance than few-shot and zero-shot learning in LLMs. Overall, Prompt Space could serve as an efficient tool for solving reasoning tasks, but also has the potential to be a few-shot learner for a wide range of applications and tasks.

ETHICS AND REPRODUCIBILITY STATEMENT

Compared to state-of-the-art methods, Prompt Space shows more competitive performance on three categories of reasoning tasks. Besides, it will significantly increase the capability and robustness of the chain-of-thought method on massive datasets. However, there are some potential limitations to consider. First, the optimal number of basis questions is observed by experimental results. Additionally, the performance of Prompt Space could be influenced by the selections of embedding models. Finally, we use an approximating method to obtain top k basis questions, which could increase the uncertainty of this method. Overall, we will continue to work on this problem to address these limitations and develop more effective and robust prompting methods.

For reproducibility, all experiments are run by gpt-35-turbo version of the public ChatGPT model from the OpenAI API with 175 billion parameters. And these baseline methods are open-sourced implementation. To aid reviewing, we summarize the statistic of ten benchmark datasets, and include configures of different embedding models and experimental settings in the supplementary materials.

REFERENCES

- Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, 2(4):433–459, 2010.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. *arXiv preprint arXiv:2210.00720*, 2022.
- Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*, 2020.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361, 2021.
- Lewis D Griffin, Bennett Kleinberg, Maximilian Mozes, Kimberly T Mai, Maria Vau, Matthew Caldwell, and Augustine Marvor-Parker. Susceptibility to influence of large language models. *arXiv preprint arXiv:2303.06074*, 2023.
- Konstantin Hebenstreit, Robert Praas, Louis P Kiesewetter, and Matthias Samwald. An automatically discovered chain-of-thought prompt generalizes to novel models and datasets. *arXiv preprint arXiv:2305.02897*, 2023.
- Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. Learning to solve arithmetic word problems with verb categorization. In *EMNLP*, pp. 523–533, 2014.

- Ellen Jiang, Kristen Olson, Edwin Toh, Alejandra Molina, Aaron Donsbach, Michael Terry, and Carrie J Cai. Promptmaker: Prompt-based prototyping with large language models. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*, pp. 1–8, 2022.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. Parsing algebraic word problems into equations. *Transactions of the Association for Computational Linguistics*, 3:585–597, 2015.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. *arXiv preprint arXiv:1705.04146*, 2017.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*, 2021.
- OpenAI. Gpt-4 technical report, 2023.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35: 27730–27744, 2022.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*, 2021.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- Laria Reynolds and Kyle McDonell. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, pp. 1–7, 2021.
- Subhro Roy, Tim Vieira, and Dan Roth. Reasoning about quantities in natural language. *Transactions of the Association for Computational Linguistics*, 3:1–13, 2015.
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*, 2021.
- Tevan Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- Timo Schick and Hinrich Schütze. It’s not just size that matters: Small language models are also few-shot learners. *arXiv preprint arXiv:2009.07118*, 2020.
- Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. Synthetic prompting: Generating chain-of-thought demonstrations for large language models. *arXiv preprint arXiv:2302.00618*, 2023.
- Grigori Sidorov, Alexander Gelbukh, Helena Gómez-Adorno, and David Pinto. Soft similarity and soft cosine measure: Similarity of features in vector space model. *Computación y Sistemas*, 18(3): 491–504, 2014.
- Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, et al. Selective annotation makes language models better few-shot learners. *arXiv preprint arXiv:2209.01975*, 2022.
- Tian-Xiang Sun, Xiang-Yang Liu, Xi-Peng Qiu, and Xuan-Jing Huang. Paradigm shift in natural language processing. *Machine Intelligence Research*, 19(3):169–183, 2022.

- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421>.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.
- Michael E Wall, Andreas Rechtsteiner, and Luis M Rocha. Singular value decomposition and principal component analysis. *A practical approach to microarray data analysis*, pp. 91–109, 2003.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*, 2022a.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in Neural Information Processing Systems*, 33:5776–5788, 2020.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022b.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- Defeng Xie, Ruichen Wang, Jian Ma, Chen Chen, Haonan Lu, Dong Yang, Fobo Shi, and Xiaodong Lin. Edit everything: A text-guided generative system for images editing. *arXiv preprint arXiv:2304.14006*, 2023.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pp. 12697–12706. PMLR, 2021.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwon Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*, 2022.

APPENDIX

The following is the supplementary Appendix for the paper. All the references are made in context of the main paper.

A DERIVATION OF THE MATRIX Q

The matrices $QQ^T \in \mathbb{R}^{m \times m}$ and $Q^TQ \in \mathbb{R}^{n \times n}$ are diagonalized, which can be rewritten as:

$$\begin{cases} QQ^T = U\Lambda_1U^T \\ Q^TQ = V\Lambda_2V^T \end{cases} \quad (\text{A.1})$$

B ALGORITHM OF PROMPT SPACE

Algorithm 1 The detailed algorithm of the proposed Prompt Space.

Input: Pre-trained text embedding models E , a question set $\mathcal{Q} = \{q_i\}_{i=1}^m$, and a large language model LLM

Parameter: A conditional variable $CoT-Zero$ for determining two cases: Prompt-Space-CoT-Zero and Prompt-Space-CoT (default $CoT-Zero=\mathbf{True}$), and the initial number of basis questions is k

Output: Answers $\{a_i\}_{i=1}^m$ from LLM

- 1: Embed each question q_i with E to yield vectors \mathbf{q}_i . Then combine all question embeddings as a matrix $Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_m]^T \in \mathbb{R}^{m \times n}$.
 - 2: Factorize the matrix Q with SVD: $Q = U\Lambda V^T$, where $Q \in \mathbb{R}^{m \times n}$.
 - 3: Find k principal components of Q , i.e., $Q_k = U_kQ$, where $U_k = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k]^T \in \mathbb{R}^{k \times m}$, and $Q_k \in \mathbb{R}^{k \times n}$.
 - 4: Calculate the similarity between basis vectors and question embeddings, and then obtain the indices of the most similar problems, i.e., $argmax(Q_k, Q^T) = [I_1, I_2, \dots, I_k]^T$, for $1 \leq I_j \leq m$.
 - 5: Construct a demonstration with the following format:
 - q_{I_1} , A: Let's think step by step. ...
 - q_{I_2} , A: Let's think step by step. ...
 - ...
 - q_{I_k} , A: Let's think step by step. ...
 - 6: **for** each question $i = 1, \dots, m$ **do**
 - 7: **if** $CoT-Zero$ **then**
 - 8: Combine the demonstration from STEP 5 with the current test question q_i :
 - ... (demonstration in STEP 5)
 - q_i
 - A: Let's think step by step.
 - 9: **else**
 - 10: Combine the demonstration from STEP 5 with the current test question q_i :
 - ... (demonstration in STEP 5)
 - q_i
 - A:
 - 11: **end if**
 - 12: Generate the output answer a_i with the prompt obtained from STEP 7 in the language model LLM
 - 13: **end for**
-

C DETAILS OF EXPERIMENTAL SETUP

C.1 DATASETS

Table C1 summarizes the basic statistics of ten benchmark datasets.

Table C1: Statistics of ten benchmark datasets.

Dataset	Answer format	# of Samples	Average words	# of basis questions	License
AddSub	Number	395	31.5	8	Unspecified
MultiArith	Number	600	31.8	8	Unspecified
SingleEq	Number	508	27.4	8	No license
AQUA-RAT	Multiple choices	254	51.9	4	Apache-2.0
SVAMP	Number	1000	31.8	8	MIT license
GSM8K	Number	1319	46.9	8	MIT license
CommonsenseQA (CSQA)	Multiple choices	1221	27.8	7	Unspecified
StrategyQA (STQA)	Yes or no	2290	9.6	6	Apache-2.0
Last Letter (Letter)	Free format	500	15.0	4	-
Coin Flip (Coin)	Yes or no	500	37.0	8	-

C.2 DEMONSTRATIONS OF BASELINES

Figure C1 shows the illustrations of five baselines, including Zero-shot, Few-shot, Manual-CoT, Zero-shot-CoT, and Auto-CoT. It is clear to see the difference between five baselines in Fig. C1.

C.3 MODELS

Table C2: Configurations of different embedding models.

Model	Configurations		
	# of Layers	Hidden size	# of Parameters
MiniLM-L6-v2	6	384	22M
E5-small	12	384	33M
E5-base	12	768	110M
E5-large	24	1024	330M
Sentence-t5-base	12	768	110M
Sentence-t5-large	24	768	336M
Sentence-t5-xl	24	768	1242M
Sentence-t5-xxl	24	768	4866M

D COMPLEMENTARY EXPERIMENTAL RESULTS

D.1 EFFECT OF BASIS QUESTIONS IN PROMPT-SPACE-COT

We additionally evaluated Prompt Space w/o CoT-Zero and Prompt-Space-CoT on ten tasks as the number of basis questions is increased from 0 to 10. The result shows that an appropriate number of basis questions is 8 for half of the arithmetic reasoning tasks (MultiArith, GSM8K, AddSub), while it varies in other arithmetic tasks. For commonsense reasoning tasks, the optimal numbers of basis questions are 6 and 9 on CSQA and STQA, respectively. Importantly, we only use one number of basis questions to achieve best performance on the Letter dataset, and that of basis questions is also small (3) on the Coin dataset. These findings indicate that our Prompt-Space-CoT needs a few basis questions to get the best performance on symbolic reasoning tasks, which dramatically reduces the cost of exemplar constructions in LLMs. For a fair comparison, we don't show the best results with optimal basis questions, while showing the results with the same number of exemplars in Table 2. Overall, these results further demonstrate that the existence of Prompt Space is significant for improving the reasoning abilities of LLMs, and reducing the cost of exemplar constructions.

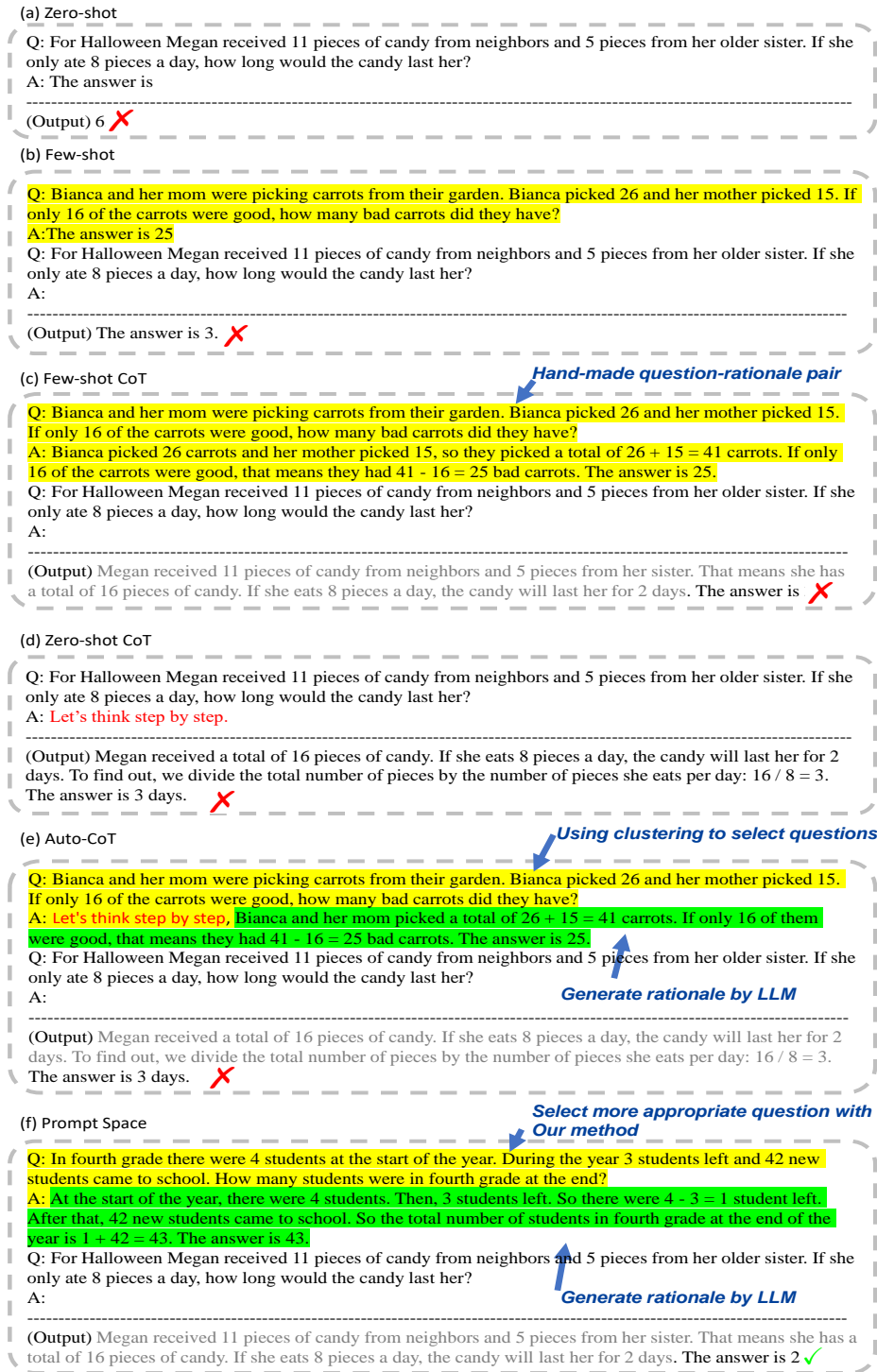


Figure C1: Demonstrations of five baselines, including Zero-shot, Few-shot, Few-shot-CoT (Manual-CoT), Zero-shot-CoT, and Auto-CoT.

Table D3: Accuracy (%) comparison of Prompt Space w/o CoT-Zero with different amounts of basis questions on ten reasoning datasets. Ten benchmark datasets contain three categories, including arithmetic reasoning, commonsense reasoning, and symbolic reasoning.

# of Basis	Arithmetic						Commonsense		Symbolic		Avg
	AddSub	MultiArith	SingleEq	AQUA-RAT	SVAMP	GSM8K	CSQA	STQA	Letter	Coin	
1	87.1	69.7	84.1	26.0	70.8	19.0	77.4	46.0	3.2	15.4	49.9
2	87.1	82.5	86.8	29.5	70.4	24.8	77.7	45.1	2.8	57.2	56.4
3	89.1	84.3	88.4	32.3	74.2	23.2	78.2	52.8	4.2	41.8	56.9
4	89.4	83.0	88.8	31.9	74.5	25.7	78.7	55.5	3.4	30.6	56.1
5	89.1	83.8	88.4	30.7	75.4	24.9	78.7	56.2	4.4	51.8	58.3
6	89.6	83.7	88.2	31.9	75.6	25.9	79.0	61.3	4.2	63.8	60.3
7	88.6	84.5	88.6	29.9	75.6	25.4	79.1	60.3	4.0	59.0	59.5
8	88.4	83.7	88.8	30.7	75.5	25.2	80.0	62.0	3.8	61.2	59.8
9	89.9	86.3	88.4	30.3	74.7	25.8	79.9	61.4	5.2	60.2	60.2
10	89.6	85.8	88.2	29.9	75.0	25.5	79.0	62.8	4.2	60.6	60.1
Best results	89.9	86.3	88.8	32.3	75.6	25.9	80.0	62.8	5.2	63.8	61.1

Table D4: Accuracy (%) comparison of Prompt-Space-CoT with different amounts of basis questions on ten reasoning datasets. Ten benchmark datasets contain three categories, including arithmetic reasoning, commonsense reasoning, and symbolic reasoning.

# of Basis	Arithmetic						Commonsense		Symbolic		Avg
	AddSub	MultiArith	SingleEq	AQUA-RAT	SVAMP	GSM8K	CSQA	STQA	Letter	Coin	
1	39.7	58.8	41.7	15.0	34.3	24.7	72.2	44.6	84.4	61.8	47.7
2	80.0	94.3	91.9	32.3	35.6	69.3	74.4	59.7	73.0	88.0	69.9
3	85.1	96.2	92.3	36.2	77.1	71.1	74.8	61.7	76.0	100.0	77.1
4	86.6	96.8	92.5	40.6	80.4	75.7	74.0	61.9	74.2	99.6	78.2
5	86.1	96.8	92.5	43.3	81.7	76.1	72.5	61.6	75.6	100.0	78.6
6	86.6	97.3	91.7	39.8	82.5	75.6	75.6	62.5	74.4	97.8	78.4
7	86.1	98.0	91.9	48.8	82.1	76.6	74.1	63.6	74.0	84.2	78.0
8	87.9	98.8	92.1	47.6	81.0	77.9	74.1	62.5	79.6	100.0	80.2
9	87.8	97.5	92.1	45.3	82.1	76.6	74.3	64.4	76.2	99.4	79.6
10	87.6	97.8	91.5	45.3	82.6	76.9	74.4	64.4	77.6	99.0	79.7
Best results	87.9	98.8	92.5	48.8	82.6	77.9	75.6	64.4	84.4	100.0	81.3

D.2 VISUALIZATION OF PROMPT SPACE

Figure D2 visualizes Prompt Space via PCA projection on ten different datasets, namely AddSub, MultiArith, SingleEq, SVAMP, AQUA-RAT, GSM8K, CSQA, STQA, Letter, and Coin. The visualization demonstrates that there exists a real space including basis vectors (questions) on each dataset. The space on each dataset may be represented by different basis questions. From Figs. D2, basis questions are located on three base planes in the three-dimensional space.

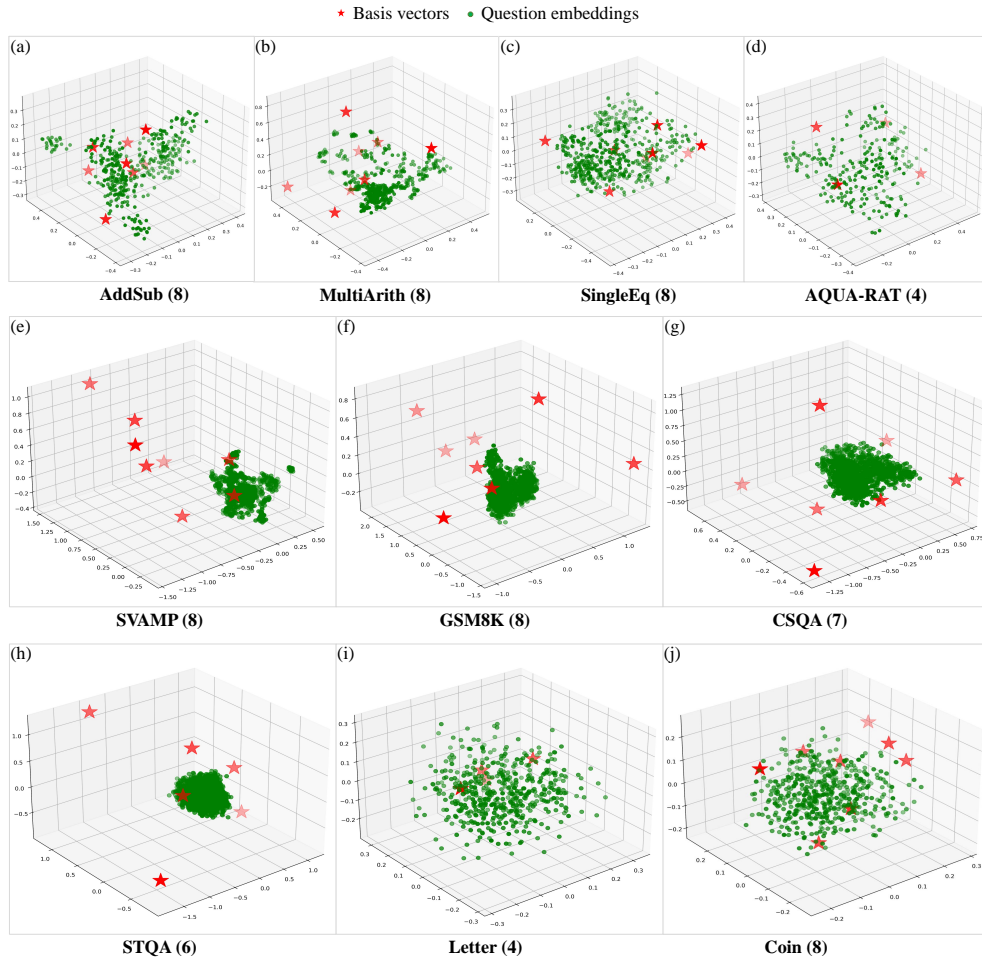


Figure D2: Visualizations of basis questions of Prompt Space on ten benchmark datasets.

E DEMONSTRATIONS OF EACH DATASETS

Table E5: Basis questions of the AddSub dataset

Q: As part of a lesson on earthquakes , a science class is researching the movement of a nearby fault line . The fault line moved 1.25 inches during the past year and 5.25 inches the year before . How far did the fault line move in all ?

A: **Let's think step by step.** The fault line moved 1.25 inches during the past year and 5.25 inches the year before. That means it moved a total of 6.5 inches in two years.

Q: There are 5 oak trees currently in the park . Park workers will plant 4 oak trees today . How many oak trees will the park have when the workers are finished ?

A: **Let's think step by step.** There are 5 oak trees currently in the park. Park workers will plant 4 oak trees today. This means that when the workers are finished, there will be a total of 9 oak trees in the park.

Q: There is 0.1666666666666666 cup of oil in Scarlett 's measuring cup . If Scarlett adds 0.6666666666666666 cup more , how much oil will be in the measuring cup ?

A: **Let's think step by step.** There is currently 0.1666666666666666 cups of oil in the measuring cup. If Scarlett adds 0.6666666666666666 more cups, that means she is adding (0.6666666666666666 + 0.1666666666666666) cups of oil, which equals 0.8333333333333334 cups of oil.

Q: Fred picked 36 limes , Alyssa picked 32 limes , and Nancy picked 35 limes and 18 pears , at the farm . How many limes were picked in total ?

A: **Let's think step by step.** First, we need to find out how many limes were picked in total. We can do this by adding up the number of limes that each person picked. $36 + 32 + 35 = 103$ So, 103 limes were picked in total.

Q: As part of a lesson on earthquakes , a science class is researching the movement of a nearby fault line . The fault line moved 1.25 inches during the past year and 5.25 inches the year before . How far did the fault line move in all ?

A: **Let's think step by step.** The fault line moved 1.25 inches during the past year and 5.25 inches the year before. That means it moved a total of 6.5 inches in two years.

Q: There are 41 short trees and 44 tall trees currently in the park . Park workers will plant 57 short trees today . How many short trees will the park have when the workers are finished ?

A: **Let's think step by step.** There are 41 short trees and 44 tall trees currently in the park. Park workers will plant 57 short trees today. This means that when the workers are finished, there will be $41 + 57 = 98$ short trees in the park.

Q: Alyssa picked 17 plums and Jason picked 10 plums . Melanie picked 35 pears . How many plums were picked in all ?

A: **Let's think step by step.** Alyssa picked 17 plums. Jason picked 10 plums. That means that, in total, 27 plums were picked.

Q: Greg and Sharon own neighboring cornfields . Greg harvested 0.4 acre of corn on Monday and Sharon harvested 0.1 acre . How many more acres did Greg harvest than Sharon ?

A: **Let's think step by step.** Greg harvested 0.4 acre of corn on Monday. That means Sharon harvested 0.1 acre of corn on Monday. That means Greg harvested 0.3 more acres of corn than Sharon.

Table E6: Basis questions of the MultiArith dataset

Q: A group of 3 friends went into a restaurant. The chef already had 6 chicken wings cooked but cooked 3 more for the group. If they each got the same amount how many would each person get?

A: **Let's think step by step.** There are a total of 9 chicken wings. Each person would get 3 chicken wings.

Q: In a video game, each enemy defeated gives you 9 points. If a level has 11 enemies total and you destroy all but 3 of them, how many points would you earn?

A: **Let's think step by step.** Each enemy defeated gives you 9 points. If a level has 11 enemies total and you destroy all but 3 of them, you would have destroyed 8 enemies. 8 times 9 equals 72 points.

Q: Wendy uploaded 45 pictures to Facebook. She put 27 pics into one album and put the rest into 9 different albums. How many pictures were in each album?

A: **Let's think step by step.** First, we know that Wendy uploaded 45 pictures in total. Second, we know that Wendy put 27 pictures into one album. That means that Wendy put the remaining 18 pictures into 9 different albums. That means that each album would have 2 pictures.

Q: Wendy was playing a video game where she scores 5 points for each treasure she finds. If she found 4 treasures on the first level and 3 on the second, what would her score be?

A: **Let's think step by step.** On the first level, Wendy scores 5 points for each treasure, so she scores 20 points in total. On the second level, Wendy scores 5 points for each treasure, so she scores 15 points in total. To find Wendy's total score, we need to add the points she scored on each level together. So we take $20 + 15$ and we get 35.

Q: Faye was playing a video game where she scores 7 points for each treasure she finds. If she found 2 treasures on the first level and 6 on the second, what would her score be?

A: **Let's think step by step.** On the first level, Faye found 2 treasures. So her score for the first level is $7 * 2 = 14$. On the second level, Faye found 6 treasures. So her score for the second level is $7 * 6 = 42$. To find her total score, we just need to add up her scores for each level. So her total score is $14 + 42 = 56$.

Q: A waiter had 9 tables he was waiting on, with 4 women and 3 men at each table. How many customers total did the waiter have?

A: **Let's think step by step.** There are 9 tables. Each table has 4 women and 3 men. That means there are 36 women and 27 men in total. Therefore, the waiter had 63 customers in total.

Q: While shopping for music online, Nancy bought 3 country albums and 5 pop albums. Each album came with a lyric sheet and had 8 songs. How many songs did Nancy buy total?

A: **Let's think step by step.** Nancy bought 3 country albums. Each album has 8 songs. So she bought 24 songs from the country albums. Nancy bought 5 pop albums. Each album has 8 songs. So she bought 40 songs from the pop albums. Nancy bought a total of $24 + 40 = 64$ songs.

Q: In a video game, each enemy defeated gives you 9 points. If a level has 11 enemies total and you destroy all but 3 of them, how many points would you earn?

A: **Let's think step by step.** Each enemy defeated gives you 9 points. If a level has 11 enemies total and you destroy all but 3 of them, you would have destroyed 8 enemies. 8 times 9 equals 72 points.

Table E7: Basis questions of the SingleEq dataset

Q: Mary, Sam, Keith, and Alyssa each have 6 baseball cards. How many baseball cards do they have in all ?

A: **Let's think step by step.** There are four people, so we need to add up all of their baseball cards. Mary has 6 baseball cards, Sam has 6 baseball cards, Keith has 6 baseball cards, and Alyssa has 6 baseball cards. $6 + 6 + 6 + 6 = 24$ So, they have 24 baseball cards in all.

Q: As part of a lesson on earthquakes, a science class is researching the movement of a nearby fault line. The fault line moved 1.25 inches during the past year and 5.25 inches the year before. How far did the fault line move in all?

A: **Let's think step by step.** The fault line moved 1.25 inches during the past year. The fault line moved 5.25 inches the year before. To find the total movement of the fault line, we need to add $1.25 + 5.25$. The fault line moved 6.5 inches in all.

Q: It rained 0.2 inches on Saturday and 0.4 inches on Sunday. How much did it rain on Saturday and Sunday combined?

A: **Let's think step by step.** The first thing we need to do is add the two numbers together. This will give us the total amount of rain that fell over the two days. $0.2 + 0.4 = 0.6$ Therefore, it rained a total of 0.6 inches on Saturday and Sunday combined.

Q: Sandy has 8 dozen books. Fred has 5 times more books than Sandy. How many books does Fred have ?

A: **Let's think step by step.** Sandy has 8 dozen books. This is 96 books. Fred has 5 times more books than Sandy. This means that Fred has 5 times 96 books, or 480 books.

Q: As part of a lesson on earthquakes, a science class is researching the movement of a nearby fault line. The fault line moved 1.25 inches during the past year and 5.25 inches the year before. How far did the fault line move in all?

A: **Let's think step by step.** The fault line moved 1.25 inches during the past year. The fault line moved 5.25 inches the year before. To find the total movement of the fault line, we need to add $1.25 + 5.25$. The fault line moved 6.5 inches in all.

Q: A DVD book holds 126 DVDs. There are 81 DVDs already in the book. How many more DVDs can be put in the book?

A: **Let's think step by step.** There are already 81 DVDs in the book. That means there are 45 DVDs that can still fit in the book.

Q: Sally bought 4 dozen eggs from the grocery store to bake some cakes. How many eggs did Sally buy ?

A: **Let's think step by step.** 1 dozen = 12 eggs 4 dozen = $4 \times 12 = 48$ eggs So Sally bought 48 eggs from the grocery store.

Q: Bert runs 2 miles every day. How many miles will Bert run in 3 weeks?

A: **Let's think step by step.** In one week, Bert will run 14 miles. In two weeks, Bert will run 28 miles. In three weeks, Bert will run 42 miles.

Table E8: Basis questions of the AQUA-RAT dataset

<p>Q: 30 is subtracted from a number, it is reduced to its one third. What is the value of 50</p> <p>A: <i>Let's think step by step.</i> First, we know that if 30 is subtracted from a number, it is reduced to one third of its original value. Next, we want to find 50% of that number. To do this, we can take one third of the number and multiply it by 2.</p> <p>Q: If Jill needed to buy 10 bottles of soda for a party in which 8 people attended, how many bottles of soda will she need to buy for a party in which 12 people are attending? Answer Choices: (A) 6 (B) 8 (C) 10 (D) 12 (E) 14</p> <p>A: <i>Let's think step by step.</i> For 8 people, Jill needs 10 bottles of soda. For 12 people, Jill needs ? bottles of soda. Since Jill needs 2.5 bottles of soda per person, she will need $2.5 * 12 = 30$ bottles of soda for 12 people. The answer is (E).</p> <p>Q: A computer routine was developed to generate two numbers (x,y) the first being a random number between 0 and 100 inclusive, and the second being less than or equal to the square root of the first. Each of the following pair satisfies the routine except Answer Choices: (A) (99,10) (B) (85,9) (C) (50,7) (D) (1,1) (E) (1,0)</p> <p>A: <i>Let's think step by step.</i> The first number, x, is a random number between 0 and 100 inclusive. The second number, y, is less than or equal to the square root of the first. That means that y can never be greater than 10. Therefore, (99,10) does not satisfy the routine.</p> <p>Q: If a man rows at the rate of 4 kmph in still water and his rate against the current is 2 kmph, then the man's rate along the current is: Answer Choices: (A) 15 kmph (B) 6 kmph (C) 12 kmph (D) 14 kmph (E) 6 kmph</p> <p>A: <i>Let's think step by step.</i> The man's rate in still water is 4 kmph. The man's rate against the current is 2 kmph. This means that the man's rate along the current must be faster than his rate in still water, but slower than his rate against the current. The answer must be between 4 kmph and 2 kmph. The answer must be 6 kmph.</p>

Table E9: Basis questions of the SVAMP dataset

Q: During summer break 61619 kids from Lawrence county stayed home and the rest went to camp. The total number of kids in Lawrence county is 91676. About how many kids in Lawrence county went to camp?

A: **Let's think step by step.** There are 91676 kids in Lawrence county. 61619 kids stayed home. That means that the rest, $91676 - 61619$, went to camp. So about 30057 kids in Lawrence county went to camp.

Q: During summer break 61619 kids from Lawrence county stayed home and the rest went to camp. The total number of kids in Lawrence county is 91676. About how many kids in Lawrence county went to camp?

A: **Let's think step by step.** There are 91676 kids in Lawrence county. 61619 kids stayed home. That means that the rest, $91676 - 61619$, went to camp. So about 30057 kids in Lawrence county went to camp.

Q: Lewis earns \$ 21 every week during the 216 weeks of harvest. If he has to pay \$ 702 tax How much money will have at the end of the harvest season?

A: **Let's think step by step.** First, let's calculate how much money he earns in total. He earns \$21 per week, and there are 216 weeks in the harvest season. So he earns $\$21 * 216 = \4536 in total. Then, we need to calculate how much tax he needs to pay. He needs to pay \$702 in tax. So the final answer is $\$4536 - \$702 = \$3834$.

Q: The grasshopper, the frog and the mouse had a jumping contest. The grasshopper jumped 25 inches. The frog jumped 18 inches farther than the grasshopper and the mouse jumped 2 inches farther than the frog. How far did the mouse jump?

A: **Let's think step by step.** The grasshopper jumped 25 inches. The frog jumped 18 inches farther than the grasshopper. This means that the frog jumped a total of $(25 + 18) = 43$ inches. The mouse jumped 2 inches farther than the frog. This means that the mouse jumped a total of $(43 + 2) = 45$ inches. Therefore, the mouse jumped 45 inches.

Q: Winter is almost here and most animals are migrating to warmer countries. There are 3 bird families living near the mountain. 26 new bird families came to live near the mountain from the arctic while 2 bird families flew away further south for winter. How many bird families were left near the mountain?

A: **Let's think step by step.** There are 3 bird families living near the mountain. 26 new bird families came to live near the mountain from the arctic. 2 bird families flew away further south for winter. $3 + 26 - 2 = 27$ bird families were left near the mountain.

Q: Danny has 12 bottle caps in his collection. He found 53 bottle caps at the park. How many bottle caps does he have now?

A: **Let's think step by step.** Danny has 12 bottle caps in his collection. He found 53 bottle caps at the park. How many bottle caps does he have now? He would have $12 + 53 = 65$ bottle caps now.

Q: The grasshopper and the frog had a jumping contest. The grasshopper jumped 13 inches. The grasshopper jumped 2 inches farther than the grasshopper. How far did the frog jump?

A: **Let's think step by step.** The grasshopper jumped 13 inches. The grasshopper jumped 2 inches farther than the grasshopper. So the frog jumped 15 inches.

Q: Lewis earns \$ 21 every week during the 216 weeks of harvest. If he has to pay \$ 702 tax How much money will have at the end of the harvest season?

A: **Let's think step by step.** First, let's calculate how much money he earns in total. He earns \$21 per week, and there are 216 weeks in the harvest season. So he earns $\$21 * 216 = \4536 in total. Then, we need to calculate how much tax he needs to pay. He needs to pay \$702 in tax. So the final answer is $\$4536 - \$702 = \$3834$.

Table E10: Basis questions of the GSM8K dataset

Q: A marketing company pays its employees on a commission-based salary system. If you sell goods worth \$1000, you earn a 30% commission. Sales over \$1000 get you an additional 10% commission. Calculate the amount of money Antonella earned if she sold goods worth \$2500.

A: **Let's think step by step.** Antonella earned a 30% commission on the first \$1000, so she earned \$300. On the remaining \$1500, she earned 10% commission. So she earned an additional \$150. In total, she earned \$450.

Q: John buys 2 pairs of shoes for each of his 3 children. They cost \$60 each. How much did he pay?

A: **Let's think step by step.** John buys 2 pairs of shoes for each of his 3 children. That means he buys 6 pairs of shoes in total. Each pair of shoes costs \$60. That means the total cost is $6 \times \$60 = \360 .

Q: A marketing company pays its employees on a commission-based salary system. If you sell goods worth \$1000, you earn a 30% commission. Sales over \$1000 get you an additional 10% commission. Calculate the amount of money Antonella earned if she sold goods worth \$2500.

A: **Let's think step by step.** Antonella earned a 30% commission on the first \$1000, so she earned \$300. On the remaining \$1500, she earned 10% commission. So she earned an additional \$150. In total, she earned \$450.

Q: Daisy bought a bag of potatoes that weighed 5 pounds. She also bought a bag of sweet potatoes that weighed 2 times as much as the potatoes and a bag of carrots that weighed 3 pounds fewer than the sweet potatoes. How many pounds of carrots did Daisy buy?

A: **Let's think step by step.** The bag of potatoes weighed 5 pounds. The bag of sweet potatoes weighed 2 times as much as the potatoes. So the sweet potatoes weighed 10 pounds. The bag of carrots weighed 3 pounds fewer than the sweet potatoes. So the carrots weighed 7 pounds.

Q: A marketing company pays its employees on a commission-based salary system. If you sell goods worth \$1000, you earn a 30% commission. Sales over \$1000 get you an additional 10% commission. Calculate the amount of money Antonella earned if she sold goods worth \$2500.

A: **Let's think step by step.** Antonella earned a 30% commission on the first \$1000, so she earned \$300. On the remaining \$1500, she earned 10% commission. So she earned an additional \$150. In total, she earned \$450.

Q: After scoring 14 points, Erin now has three times more points than Sara, who scored 8. How many points did Erin have before?

A: **Let's think step by step.** Erin now has 3 times more points than Sara, who scored 8. That means that if we take away 8 points from Erin, she will still have 3 times more points than Sara. That means that Erin had 24 points before.

Q: Dijana and Anis live near a lake, and every weekend they go out rowing into the lake. On a Sunday morning, both went out rowing, and Dijana rowed for 50 miles the whole day. Anis rowed $\frac{1}{5}$ times more miles than Dijana. Calculate the total distance the two of them rowed on that day.

A: **Let's think step by step.** Dijana rowed 50 miles, so Anis rowed $\frac{1}{5}$ more than that, which is 10 miles. So the total distance they rowed is 60 miles.

Q: Kyle bought last year's best-selling book for \$19.50. This is with a 25% discount from the original price. What was the original price of the book?

A: **Let's think step by step.** The original price of the book was \$19.50 with a 25% discount. This means that the original price was \$26. The original price was \$26.

Table E11: Basis questions of the CSQA dataset

Q: Where can someone get a new saw? Answer Choices: (A) hardware store (B) toolbox (C) logging camp (D) tool kit (E) auger

A: **Let's think step by step.** Where would you go to buy a new saw? A hardware store would be a good place to start.

Q: The man was going fishing instead of work, what is he seeking? Answer Choices: (A) food (B) relaxation (C) killing (D) missing morning cartoons (E) boredom

A: **Let's think step by step.** The man is going fishing instead of work. This means that the man would rather go fishing than work. The man is seeking relaxation.

Q: What do animals do when an enemy is approaching? Answer Choices: (A) feel pleasure (B) procreate (C) pass water (D) listen to each other (E) sing

A: **Let's think step by step.** When an animal notices an enemy approaching, its first instinct is usually to flee. If fleeing is not an option, the animal may try to hide. If hiding is not an option, the animal may try to defend itself.

Q: I want to buy a gong, where should I look for one? Answer Choices: (A) orchestra (B) church (C) chinese temple (D) chinatown (E) music store

A: **Let's think step by step.** A gong is a percussion instrument, so we can eliminate (A) and (B). It is also likely that a gong would be used in a Chinese temple, so (C) and (D) are good choices. However, the best answer is probably (E), since a music store would have a wider selection of gongs to choose from.

Q: I want to buy a gong, where should I look for one? Answer Choices: (A) orchestra (B) church (C) chinese temple (D) chinatown (E) music store

A: **Let's think step by step.** A gong is a percussion instrument, so we can eliminate (A) and (B). It is also likely that a gong would be used in a Chinese temple, so (C) and (D) are good choices. However, the best answer is probably (E), since a music store would have a wider selection of gongs to choose from.

Q: What is it called when a person with mental illness is able to lead a relatively normal life? Answer Choices: (A) managed (B) effectively treated (C) recur (D) cause delusion (E) illusion

A: **Let's think step by step.** The first thing to do is to identify the key words in the question. In this case, the key words are "mental illness," "relatively normal life," and "able to lead." Next, we need to think about what these key words mean. "Mental illness" refers to a condition that affects a person's thoughts, emotions, and behavior. "Relatively normal life" means that the person is able to function in society, even though they may have some symptoms of their mental illness. "Able to lead" means that the person is able to take care of themselves and live independently.

Q: What is the feeling of one having fun? Answer Choices: (A) smiling (B) pleasure (C) hurt (D) injuries (E) laughter

A: **Let's think step by step.** The feeling of one having fun is generally associated with positive emotions, such as happiness, laughter, and pleasure. Therefore, we can eliminate choices (C) and (D). Next, we can look at the remaining choices and try to determine which one is the best fit. Both choices (A) and (E) could work, but (E) is a more specific answer, so it is probably the correct answer.

Table E12: Basis questions of the STQA dataset

Q: Is Bern located east of Paris?

A: **Let's think step by step.** First, we need to know the coordinates of Paris and Bern. According to Google Maps, Paris is located at 48.856614, 2.3522219000000177 and Bern is located at 46.947922, 7.4446080000000255. Next, we need to compare the longitude of Paris (2.3522219000000177) with the longitude of Bern (7.4446080000000255). Since the longitude of Bern is greater than the longitude of Paris, Bern is located east of Paris.

Q: If he were poor, would Christopher Reeve have lived?

A: **Let's think step by step.** If he were poor, he would not have had the money to get the best medical care possible. If he did not have the best medical care possible, his injuries would not have been treated as well as they were. If his injuries were not treated as well as they were, he would not have lived as long as he did. So, in short, if he were poor, he would not have lived as long as he did.

Q: Would it be unusual to see frost in September in Texas?

A: **Let's think step by step.** First, what is the definition of frost? Frost is defined as a thin layer of ice that forms on surfaces when the temperature outside is below freezing. In Texas, the average temperature in September is around 80 degrees Fahrenheit. So, no, it would not be unusual to see frost in September in Texas.

Q: Can an American black bear swallow a sun bear whole?

A: **Let's think step by step.** The average adult American black bear weighs between 250 and 550 pounds. The average adult sun bear weighs between 60 and 80 pounds. So an American black bear could theoretically swallow a sun bear whole.

Q: Is CAS number 8009-03-8 harmful for a rash?

A: **Let's think step by step.** CAS number 8009-03-8 is the CAS number for Shea Butter. Shea Butter is a natural fat that is extracted from the nut of the African shea tree. It is solid at room temperature but melts when it comes into contact with the skin. There is no evidence that Shea Butter is harmful for a rash. In fact, it is often used as a natural remedy for skin conditions such as eczema and psoriasis.

Q: If Martin Luther did one theses a day would he run out in half a year?

A: **Let's think step by step.** There are 365 days in a year. If Martin Luther did one theses a day, he would do 365 theses in a year. Therefore, he would not run out in half a year.

Table E13: Basis questions of the Letter dataset

Q: Take the last letters of each words in "Ever Gio Elia Ramesh" and concatenate them.

A: **Let's think step by step.** The last letter of "Ever" is "r". The last letter of "Gio" is "o". The last letter of "Elia" is "a". The last letter of "Ramesh" is "h". So, the final answer is "roah".

Q: Take the last letters of each words in "Marian Joanne Darrin Rohit" and concatenate them.

A: **Let's think step by step.** The last letter of "Marian" is "n". The last letter of "Joanne" is "e". The last letter of "Darrin" is "n". The last letter of "Rohit" is "t". So the final answer would be "nent".

Q: Take the last letters of each words in "Craig Dillon Troy Griselda" and concatenate them.

A: **Let's think step by step.** The last letter of the first word is "g". The last letter of the second word is "y". The last letter of the third word is "a". The last letter of the fourth word is "a". Therefore, the final answer is "gyaa".

Q: Take the last letters of each words in "Liliana Quincy Bart Makayla" and concatenate them.

A: **Let's think step by step.** The last letter of "Liliana" is "a". The last letter of "Quincy" is "y". The last letter of "Bart" is "t". The last letter of "Makayla" is "a". So the final answer is "ayta".

Table E14: Basis questions of the Coin Dataset

<p>Q: A coin is heads up. Miguel Angel flips the coin. Saul does not flip the coin. Brady does not flip the coin. Darryl flips the coin. Is the coin still heads up? Note that "flip" here means "reverse".</p> <p>A: Let's think step by step. Miguel flips the coin, so it's now tails up. Saul doesn't flip the coin, so it remains tails up. Brady doesn't flip the coin, so it remains tails up. Darryl flips the coin, so it's now heads up again.</p>
<p>Q: A coin is heads up. Jacky flips the coin. Socorro flips the coin. Mark flips the coin. Wanda does not flip the coin. Is the coin still heads up? Note that "flip" here means "reverse".</p> <p>A: Let's think step by step. Jacky flips the coin, so it's now tails up. Socorro flips the coin, so it's now heads up again. Mark flips the coin, so it's now tails up again. Wanda does not flip the coin, so it remains tails up.</p>
<p>Q: A coin is heads up. Claudia flips the coin. Cole does not flip the coin. Matthew does not flip the coin. Juan Pablo does not flip the coin. Is the coin still heads up? Note that "flip" here means "reverse".</p> <p>A: Let's think step by step. Claudia flips the coin, so it is now tails up. Cole does not flip the coin, so it remains tails up. Matthew does not flip the coin, so it remains tails up. Juan Pablo does not flip the coin, so it remains tails up. So the answer is no, the coin is not heads up.</p>
<p>Q: A coin is heads up. Teresa flips the coin. Reid does not flip the coin. Karin does not flip the coin. Gracie does not flip the coin. Is the coin still heads up? Note that "flip" here means "reverse".</p> <p>A: Let's think step by step. If Teresa flips the coin, then it is not heads up. If Reid flips the coin, then it is not heads up. If Karin flips the coin, then it is not heads up. If Gracie flips the coin, then it is not heads up. So, the coin is not heads up.</p>
<p>Q: A coin is heads up. Miguel Angel flips the coin. Saul does not flip the coin. Brady does not flip the coin. Darryl flips the coin. Is the coin still heads up? Note that "flip" here means "reverse".</p> <p>A: Let's think step by step. Miguel flips the coin, so it's now tails up. Saul doesn't flip the coin, so it remains tails up. Brady doesn't flip the coin, so it remains tails up. Darryl flips the coin, so it's now heads up again.</p>
<p>Q: A coin is heads up. Aj does not flip the coin. Jd flips the coin. Maddie does not flip the coin. Francisca does not flip the coin. Is the coin still heads up? Note that "flip" here means "reverse".</p> <p>A: Let's think step by step. Aj does not flip the coin, so the coin remains heads up. Jd flips the coin, so the coin is now tails up. Maddie does not flip the coin, so the coin remains tails up. Francisca does not flip the coin, so the coin remains tails up. So, the final answer is that the coin is tails up.</p>
<p>Q: A coin is heads up. Tomas does not flip the coin. Nic flips the coin. Zoila flips the coin. Calvin does not flip the coin. Is the coin still heads up? Note that "flip" here means "reverse".</p> <p>A: Let's think step by step. Tomas does not flip the coin, so the coin remains heads up. Nic flips the coin, so the coin is now tails up. Zoila flips the coin, so the coin is now heads up again. Calvin does not flip the coin, so the coin remains heads up.</p>
<p>Q: A coin is heads up. Tere does not flip the coin. Niko flips the coin. Keith does not flip the coin. Conner does not flip the coin. Is the coin still heads up? Note that "flip" here means "reverse".</p> <p>A: Let's think step by step. Tere does not flip the coin, so the coin remains heads up. Niko flips the coin, so the coin is now tails up. Keith does not flip the coin, so the coin remains tails up. Conner does not flip the coin, so the coin remains tails up. So, the final answer is that the coin is tails up.</p>
